

# Fingerprinting the Checker Policies of Parallel File Systems

Runzhou Han, Duo Zhang, Mai Zheng

IOWA STATE UNIVERSITY

# Parallel File Systems (PFSes)

- PFS is the cornerstone of high performance computing
  - Optimized for highly concurrent access



#### PFS Failures: Real-World Cases

Thursday, September 1, 2016 at 9:10:12 AM Central Daylight Time

Subject: Update: HPCC Power Outage

Date: Monday, January 11, 2016 at 8:50:17 AM Central Standard Time

From: HPCC - Support

Attachments: image001.png, image003.png

TEXAS TECH UNIVERSITY Information Technology Division

High Performance Computing Center

To All HPCC Customers and Partners,

As we have informed you earlier, the Experimental Sciences Building experienced a major power outage Sunday, Jan. 3 and another set of outages Tuesday, Jan. 5 that occurred while file systems were being recovered from the first outage. As a result, there were major losses of important parts of the file systems for the work, scratch and certain experimental group special Lustre areas.

The HPCC staff have been working continuously since these events on recovery procedures to try to restore as much as possible of the affected file systems. These procedures are extremely time-consuming, taking days to complete in some cases. Although about a third of the affected file systems have been recovered, work continues on this effort and no time estimate is possible at present.

User home areas have been recovered successfully. At present, no user logins are being permitted while recovery efforts proceed on the remaining Lustre areas. Your understanding and patience are appreciated.

If you have questions, please contact us at hpccsupport@ttu.edu or 806-742-4350. Thanks.

Sincerely, HPCC Staff



#### PFS Failures: Real-World Cases

Thursday, September 1, 2016 at 9:10:12 AM Central Daylight Time

Subject: Update: HPCC Power Outage

Date: Monday, January 11, 2016 at 8:50:17 AM Central Standard Time

From: HPCC - Support

Attachments: image001.png, image003.png

Information Technology Division

High Performance Computing Center

To All HPCC Customers and Partners,

As we have informed you earlier, the Experimental Sciences Building experienced a major power outage Sunday, Jan. 3 and another set of outages Tuesday, Jan. 5 that occurred while file systems were being recovered from the first outage. As a result, there were major losses of important parts of the file systems for the work, scratch and certain experimental group special Lustre areas.

The HPCC staff have been working continuously since these events on recovery procedures to try to restore as much as possible of the affected file systems. These procedures are extremely time-consuming, taking days to complete in some cases. Although about a third of the affected file systems have been recovered, work continues on this effort and no time estimate is possible at present.

User home areas have been recovered successfully. At present, no user logins are being permitted while recovery efforts proceed on the remaining Lustre areas. Your understanding and patience are appreciated.

If you have questions, please contact us at <u>hpccsupport@ttu.edu</u> or 806-742-4350. Thanks.

Sincerely, HPCC Staff

Case1: HPCC Power Outage

The graph below offers a vivid example of how bad things got at a prestigious American university<sup>2</sup> which suffered frequent HPC storage outages and took several days to get their systems back up and running. Recovery experienced at this university<sup>3</sup> is shown in the graph below with additional detail available via the link in the footnote. It shows an outage that started on a Monday and wasn't fully recovered from until Sunday.



#### Case2: ACCRE Storage Outage\*

\* Hyperion Research survey of HPC organizations done for Panasas

#### PFS Failures: More Frequent/Expensive Than You Thought

#### **Some statistics\***:

The average HPC storage system failure frequency is **9.8 failures**/year

#### ≈half

of HPC sites experience storage system failures 1/month or more frequently

\* Hyperion Research survey of HPC organizations done for Panasas

#### PFS Failures: More Frequent/Expensive Than You Thought

#### Some statistics\*:

The average HPC storage system failure frequency is 9.8 failures/year

Downtime ranges from  $1 \, day_{\downarrow}$  to  $1 \, week_{\uparrow}$ 

#### ≈half

of HPC sites experience storage system failures 1/month or more frequently

#### 40%

of HPC sites typically took more than 2 weeks to restore their storage systems

\* Hyperion Research survey of HPC organizations done for Panasas

#### PFS Failures: More Frequent/Expensive Than You Thought

#### Some statistics\*:

The average HPC storage system failure frequency is 9.8 failures/year

Downtime ranges from  $1 \, day_{\downarrow} \text{ to } 1 \, week_{\uparrow}$ 

A single day of downtime costs from  $\$100K\downarrow \text{ to }\$1M\uparrow$ 

#### ≈half

of HPC sites experience storage system failures 1/month or more frequently 40% of HPC sites typically took more than 2 weeks to restore their storage systems

Average downtime cost is \$127K/day

#### • Typical PFS architecture



• Typical PFS architecture



- Many PFSes are designed with a *checker* component
  - e.g., LFSCK for Lustre, BeeGFS-FSCK for BeeGFS, PV2FS-FSCK for OrangeFS



**BeeGFS**<sup>®</sup>



• Typical PFS architecture



- Many PFSes are designed with a *checker* component
  - e.g., LFSCK for Lustre, BeeGFS-FSCK for BeeGFS, PV2FS-FSCK for OrangeFS
  - Detect and repair inconsistencies Lustre<sup>\*</sup>





• Typical PFS architecture



- Many PFSes are designed with a *checker* component
  - e.g., LFSCK for Lustre, BeeGFS-FSCK for BeeGFS, PV2FS-FSCK for OrangeFS

BeeGFS<sup>®</sup>

- Detect and repair inconsistencies Lustre<sup>\*</sup>
- FSCKs have predefined *checker policies*

**Örange**FS

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID
  - OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID
  - OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID
  - OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID
  - OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID
  - OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID
  - OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  Corruption 1 MDT-object's LOV EA matches to OST-object's FID
  - OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

Lustre's LFSCK Policy: mapping between MDT-object and OST-object

Corruption 1 • MDT-object's LOV EA matches to OST-object's FID Fixed!

• OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID

Corruption 2 • OST-object's Parent FID matches to MDT-object's FID



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID

Corruption 2 • OST-object's Parent FID matches to MDT-object's FID Cannot be fixed!



Structures	Meaning	
xattr	inode extended attribute	
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

- Lustre's LFSCK Policy: mapping between MDT-object and OST-object
  - MDT-object's LOV EA matches to OST-object's FID

Corruption 2 • OST-object's Parent FID matches to MDT-object's FID Cannot be fixed!

#### LFSCK's policy is incomplete!



λατι		
FID	a global ID of an Lustre object	
LOV EA	stores child object's FID	
Parent FID	stores parent object's FID	

#### Our Contributions

- A systematic approach to analyze PFS checker policies
  - PFS type-aware fault injection
  - PFS consistency model & taxonomy

## Our Contributions

- A systematic approach to analyze PFS checker policies
  - PFS type-aware fault injection
  - PFS consistency model & taxonomy
- A comprehensive study on the checkers of two widely used PFSes



- Has exposed 33 suboptimal repairs
- Has exposed 2 abnormal behaviors (e.g., kernel panic), which has led to 1 new patch on Lustre





# Outline

- Motivation & Contributions
- Methodology
  - PFS type-aware fault injection
  - Fault models
  - PFS consistency model
  - PFS checker taxonomy
- Experimental Result
- Conclusion & Future Work

# PFS Type-aware Fault Injection

- Why type-aware fault injection
  - Key observation
    - PFS metadata and the local file system metadata are closely correlated
    - E.g., Both Lustre and BeeGFS has metadata structures stored in inode extended attribute (xattr)

PFS name	metadata structures in xattr	shared metadata structures with Ext4 inode
Lustre	FID, LOV EA, parent FID, linkEA	nlink
BeeGFS	fhgfs	nlink, size

- Benefits of fine-grained fault injection
  - reveal PFS checker policies precisely
  - Enable analyzing the contract between PFS checker and local FS

#### Fault Models

- Four fault models to capture the typical corruptions that may occur in the local storage stack and be exposed to the PFS checker
  - *junk, out-of-sync, zero, duplicate*

### Fault Models

- Four fault models to capture the typical corruptions that may occur in the local storage stack and be exposed to the PFS checker
  - junk, out-of-sync, zero, duplicate
- Fault model #1: junk
  - Bytes of the on-disk structure are replaced by random values
    - Caused by disk corruptions, local FS bugs, etc



### Fault Models

- Fault model #2: *out-of-sync* 
  - In-memory copy of the structure is inconsistent with on-disk copy
    - Caused by software bugs ,memory/disk corruptions, etc



• Please refer to our paper for fault models #3 & #4

- General principles that PFS checkers should ensure to maintain PFS integrity
  - Applicable to diverse PFSes
  - Include the definition of Consistency Group & 6 consistency rules

- General principles that PFS checkers should ensure to maintain PFS integrity
  - Applicable to diverse PFSes
  - Include the definition of Consistency Group & 6 consistency rules
- Consistency Group (CG)
  - Include an MDT-object and all its associated child OST-objects



- General principles that PFS checkers should ensure to maintain PFS integrity
  - Applicable to diverse PFSes
  - Include the definition of Consistency Group & 6 consistency rules
- Consistency Group (CG)
  - Include an MDT-object and all its associated child OST-objects
  - Consistency rules
    - CG-rule1: every object in a CG should be consistent individually



- General principles that PFS checkers should ensure to maintain PFS integrity
  - Applicable to diverse PFSes
  - Include the definition of Consistency Group & 6 consistency rules
- Consistency Group (CG)
  - Include an MDT-object and all its associated child OST-objects
  - Consistency rules
    - *CG-rule2*: one MDT-object of a client directory maps to no child OST-object



- General principles that PFS checkers should ensure to maintain PFS integrity
  - Applicable to diverse PFSes
  - Include the definition of Consistency Group & 6 consistency rules
- Consistency Group (CG)
  - Include an MDT-object and all its associated child OST-objects
  - Consistency rules
    - *CG-rule3*: one MDT-object of a client file maps to at least one child OST-object



- General principles that PFS checkers should ensure to maintain PFS integrity
  - Applicable to diverse PFSes
  - Include the definition of Consistency Group & 6 consistency rules
- Consistency Group (CG)
  - Include an MDT-object and all its associated child OST-objects
  - Consistency rules
    - *CG-rule4*: one OST-object maps to one and only one parent MDT-object



- General principles that PFS checkers should ensure to maintain PFS integrity
  - Applicable to diverse PFSes
  - Include the definition of Consistency Group & 6 consistency rules
- Consistency Group (CG)
  - Include an MDT-object and all its associated child OST-objects
  - Consistency rule
    - CG-rule5: the mapping b/w a parent MDT-object and a child OST-object is bidirectional



- General principles that PFS checkers should ensure to maintain PFS integrity
  - Applicable to diverse PFSes
  - Include the definition of Consistency Group & 6 consistency rules
- Consistency Group (CG)
  - Include an MDT-object and all its associated child OST-objects
  - Consistency rules
    - CG-rule6: an object violating previous rules may only exist in a specified location



- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

- A general characterization of checker policies
  - Qualitatively measures the policies
  - Enable cross-PFS comparison
- Include 4 Detection levels & 4 Repair levels based on consistency model

Detection levels	Definition	Repair levels	Definition
D <sub>abn.</sub>	PFS checker behaves abnormally w/o reporting detection results	R <sub>wro.</sub>	PFS checker fixes CG corruptions in a wrong way
D <sub>zero</sub>	PFS checker finishes normally but misses all CG corruptions	R <sub>zero</sub>	PFS checker reports failure on repair
D <sub>par.</sub>	PFS checker partially detects CG corruptions	R <sub>par.</sub>	PFS checker partially fixes CG corruptions
D <sub>com.</sub>	PFS checker detects CG corruptions completely	R <sub>com.</sub>	PFS checker fixes corruptions and CGs're valid again

## Outline

- Motivation & Contributions
- Methodology
  - PFS type-aware fault injection
  - Fault models
  - PFS consistency model
  - PFS checker taxonomy
- Experimental Results
- Conclusion & Future Work

Lustre Structures	ju	nk	ze	ro	dupli	duplicate      out-of-sync      BeeGFS Structures      junk		ze	zero		duplicate		f-sync				
MDT-object	_	_	-	_	_	_	D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-name (MDT-object)	_	-			_		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_	_	D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-ID (MDT-object)	_	_		_	-		D <sub>com.</sub>	R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	_	D <sub>zero</sub>	R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub> R <sub>zero</sub>		D <sub>zero</sub>	R <sub>zero</sub>	-	-	D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	$R_{zero}$	D <sub>com.</sub>	R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	content directory	—		-			_	D <sub>com.</sub>	R <sub>com.</sub>
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_		D <sub>zero</sub>	R <sub>zero</sub>	nlink	*D <sub>com.</sub> R <sub>zero</sub>		*D <sub>com.</sub>	R <sub>zero</sub>	_			
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	_	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub> R <sub>zero</sub> *		*D <sub>com.</sub>	R <sub>zero</sub>	_		—	
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	D <sub>zero</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	_	_	_								47	

Lustre Structures	ju	nk	zero		duplicate out-of-sync		f-sync	BeeGFS Structures	junk		zero		duplicate		out-of-sync		
MDT-object	-	_	_		—		D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-name (MDT-object)	—		-		- –		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>		-	D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-ID (MDT-object)	-	-	_	-		-		R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	-	D <sub>zero</sub>	R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub> R <sub>zero</sub>		D <sub>zero</sub>	R <sub>zero</sub>	_	-	D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	$R_{zero}$	D <sub>com.</sub>	R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	content directory	—		-	_	_		D <sub>com.</sub>	R <sub>com.</sub>
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_		D <sub>zero</sub>	R <sub>zero</sub>	nlink	*D <sub>com.</sub> R <sub>zero</sub>		*D <sub>com.</sub>	R <sub>zero</sub>	_			
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	-	-	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub> R <sub>zero</sub>		*D <sub>com.</sub>	R <sub>zero</sub>	-	-	_	
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	D <sub>zero</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	_	-								48	

Lustre Structures	ju	nk	ze	ro	dupli	icate out-of-sync		f-sync	BeeGFS Structures	junk		zero		duplicate		out-of-sync	
MDT-object	_	_	-	_	_	-	D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-name (MDT-object)	_	—		-	_		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_	-	D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-ID (MDT-object)	_	_		-	_		D <sub>com.</sub>	R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>		-	D <sub>zero</sub>	R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	-	D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	$R_{zero}$	D <sub>com.</sub>	R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	content directory	—		-	-	_		D <sub>com.</sub>	R <sub>com.</sub>
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	—		D <sub>zero</sub>	R <sub>zero</sub>	nlink	*D <sub>com.</sub> R <sub>zero</sub>		*D <sub>com.</sub>	R <sub>zero</sub>	_			
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub> R <sub>zero</sub>		*D <sub>com.</sub>	R <sub>zero</sub>	_	-	—	
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	D <sub>zero</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	_	-								49	

Lustre Structures	ju	nk	ze	ro	duplicate		out-of-sync		BeeGFS Structures	junk		zero		duplicate		out-of-sync	
MDT-object	-	_	-	_	—		D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-name (MDT-object)	—		-		_		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	_		R <sub>zero</sub>	dentry-by-ID (MDT-object)	-		-		_		D <sub>com.</sub>	R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	—		D <sub>zero</sub>	R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_	-	D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	D <sub>com.</sub> R <sub>zero</sub>		R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	content directory	-	_	-	-	_		D <sub>com.</sub>	R <sub>com.</sub>
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_			R <sub>zero</sub>	nlink	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	_	_
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	_	_
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	D <sub>zero</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_			_								50	

Lustre Structures	ju	nk	ze	ro	duplicate		out-of-sync		BeeGFS Structures	junk		zero		duplicate		out-of-sync	
MDT-object	_	-	-	_	-	_		R <sub>zero</sub>	dentry-by-name (MDT-object)	_		_		_		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>		_		R <sub>zero</sub>	dentry-by-ID (MDT-object)	—		-			-	D <sub>com.</sub>	R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	_		R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub> R <sub>zero</sub> D		D <sub>zero</sub>	R <sub>zero</sub>	_		D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	D <sub>com.</sub> R <sub>zero</sub>		R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	content directory	_	_	-	_	_		D <sub>com.</sub>	R <sub>com.</sub>
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_	_		R <sub>zero</sub>	nlink	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub> R <sub>zero</sub>					
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub> R <sub>zero</sub> *D <sub>com.</sub> R <sub>zero</sub>		R <sub>zero</sub>	-		_		
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	D <sub>zero</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	_	_	-								51	

• 14 cases: checkers repair CG corruptions completely

Lustre Structures	ju	nk	ze	ro	dupli	duplicate out-		f-sync	BeeGFS Structures	junk		zero		duplicate		out-of-sync	
MDT-object	-	_	_	_		_		R <sub>zero</sub>	dentry-by-name (MDT-object)	—		_		_		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_	_		R <sub>zero</sub>	dentry-by-ID (MDT-object)	—		-		_		D <sub>com.</sub>	R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	_		R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	_	D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	D <sub>com.</sub> R <sub>zero</sub>		R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub> R <sub>zero</sub>		D <sub>zero</sub>	R <sub>zero</sub>	content directory	_		_		_		D <sub>com.</sub>	R <sub>com.</sub>
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_	_		R <sub>zero</sub>	nlink	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	-	_	-	_
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	_	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	_	_
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	D <sub>zero</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	_		_								52	

• 18 cases: checkers detects CG corruptions but can't repair completely

Lustre Structures	ju	nk	zero		duplicate		out-of-sync		BeeGFS Structures	junk		zero		duplicate		out-of-sync	
MDT-object	-	_	—		-		D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-name (MDT-object)	—		_		_		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	_		R <sub>zero</sub>	dentry-by-ID (MDT-object)	—		—		_		D <sub>com.</sub>	R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	—		D <sub>zero</sub>	R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_		D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	content directory	-	_	_	_			D <sub>com.</sub>	R <sub>com.</sub>
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_		D <sub>zero</sub>	R <sub>zero</sub>	nlink	*D <sub>com.</sub>	*D <sub>com.</sub> R <sub>zero</sub> *D <sub>com.</sub> R <sub>zero</sub>		R <sub>zero</sub>	_		_	
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	-	_	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_		-	-
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	Dzero	R <sub>zero</sub>	*D <sub>com</sub>	R <sub>zero</sub>	_		_	_								53	

- 12 cases: checkers only check the in-memory copy of the structure
  - Could potentially miss corruptions of on-disk structures

Lustre Structures	ju	nk	zero		duplicate		out-of-sync		BeeGFS Structures	junk		zero		duplicate		out-of-sync	
MDT-object	-	_	-	_	_	_		R <sub>zero</sub>	dentry-by-name (MDT-object)	—		-		_		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	—		D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-ID (MDT-object)	-		_		_		D <sub>com.</sub>	R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	—		D <sub>zero</sub>	R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	-	_	D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	D <sub>com.</sub> R <sub>zero</sub>		R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub> R <sub>zero</sub>		D <sub>zero</sub>	R <sub>zero</sub>	content directory	-	_					D <sub>com.</sub>	R <sub>com.</sub>
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	—		D <sub>zero</sub>	R <sub>zero</sub>	nlink	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>				
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	-	_	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	-	_	_	_
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	Dzero	R <sub>zero</sub>	*D <sub>com</sub>	Rzero	_	_		_	-							54	

#### • 2 cases: LFSCK triggers kernel panic

- Has been confirmed by developers and led to 1 new patch
  - WhamCloud Community Jira: <u>LU-13980</u>, 09/26/2020

Lustre Structures	ju	nk	zero		duplicate		out-of-sync		BeeGFS Structures	junk		zero		duplicate		out-of-sync	
MDT-object	_	_	_				D <sub>abn.</sub>	R <sub>zero</sub>	dentry-by-name (MDT-object)	—				_		D <sub>com.</sub>	R <sub>com.</sub>
OST-object	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_			R <sub>zero</sub>	dentry-by-ID (MDT-object)	—		_	_	_		D <sub>com.</sub>	R <sub>com.</sub>
llog record	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	—		D <sub>zero</sub>	R <sub>zero</sub>	chunk (OST-object)	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>		-	D <sub>par.</sub>	R <sub>zero</sub>
FID on MDT	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	fghfs	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	R <sub>wro.</sub>	D <sub>com.</sub>	D <sub>com.</sub> R <sub>zero</sub>		R <sub>wro.</sub>
FID on OST	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	R <sub>zero</sub>	D <sub>com.</sub>	D <sub>com.</sub> R <sub>zero</sub> D <sub>z</sub>		R <sub>zero</sub>	content directory			_			D <sub>com.</sub>	R <sub>com.</sub>	
FLDB	D <sub>zero</sub>	R <sub>zero</sub>	D <sub>zero</sub>	R <sub>zero</sub>	_	_		R <sub>zero</sub>	nlink	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	-	-		_
OI table	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	_	*D <sub>com.</sub>	R <sub>com.</sub>	size	*D <sub>com.</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_	-	_	_
LOV EA	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>	*D <sub>com.</sub>	R <sub>com.</sub>									
PFID	D <sub>par.</sub>	R <sub>par.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>zero</sub>	R <sub>zero</sub>									
linkEA	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>	D <sub>com.</sub>	+R <sub>com.</sub>	D <sub>com.</sub>	R <sub>com.</sub>									
nlink	D <sub>zero</sub>	R <sub>zero</sub>	*D <sub>com.</sub>	R <sub>zero</sub>	_		_	-								55	

## Outline

- Motivation & Contributions
- Methodology
  - PFS type-aware fault injection
  - Fault models
  - PFS consistency model
  - PFS checker taxonomy
- Experimental Result
- Conclusion & Future Work

### Conclusion & Future Work

- A systematic approach to study PFS checkers
  - Has led to a new patch on Lustre
- Future work
  - More automation (e.g., apply fuzzing techniques)
  - Study other PFSes (e.g., OrangeFS, Ceph)
  - Improve PFS checkers
    - Policy completeness
    - Performance

### Conclusion & Future Work

- A systematic approach to study PFS checkers
  - Has led to a new patch on Lustre
- Future work
  - More automation (e.g., apply fuzzing techniques)
  - Study other PFSes (e.g., OrangeFS, Ceph)
  - Improve PFS checkers
    - Policy completeness
    - Performance

# Thank you & Questions?